

[Ve]tting [R]esponse [I]ntegrity from cross-[T]alk dependencies in [A]dversarial [S]urveys: *“Who can catch a liar?”*

Robert Gibbons, Ishanu Chattopadhyay^{1,4,5,7}★

¹Department of Medicine, University of Chicago, Chicago, IL 60637, USA

⁴Committee on Quantitative Methods in Social, Behavioral, and Health Sciences, University of Chicago, Chicago, IL 60637, USA

⁵Committee on Genetics, Genomics & Systems Biology, University of Chicago, Chicago, IL 60637, USA

★To whom correspondence should be addressed: e-mail: ishanu@uchicago.edu.

Abstract: Faking symptoms of a mental health disorder or *malinger* can confound diagnostic interviews and hinder clinical psychiatric assessment. We introduce an artificial intelligence (AI) framework for detecting malingering in Post-Traumatic Stress Disorder (PTSD) diagnoses. Traditional psychiatric diagnoses methods, predominantly structured interviews, are susceptible to fabrication, which is especially problematic in the context of PTSD due to subjective stressors, presentation similarity, and potential of secondary financial gain from positive diagnoses. Algorithm *veRITAS* employs a generative AI to infer statistical dependencies inherent in true response patterns, and flags responses which violate such subtle constraints. These emergent response patterns are challenging to mimic even with psychiatric training, providing a robust mechanism for recognizing fabrications. With a study sample of $n = 624$ patients, *veRITAS* is estimated to have a minimum Area Under the Curve (AUC) of 95.0%, with one sensitivity and specificity pair of 92% and 95% respectively. Thus our tool offers a objective, disease-specific, rapid (≤ 2 min) approach to detect fake PTSD, and if adopted, can ensure that healthcare resources reach those genuinely in need, and maintain the integrity of clinical data. Moreover, the ability to diagnose patients who might be malingering due to other mental health conditions, poverty or socioeconomic compulsions can substantially improve outcomes in low SES communities.

INTRODUCTION

Diagnosis of mental health disorders typically rely on structured interviews¹ susceptible to intentional fabrication of symptoms^{2,3}, often referred to as “malingering”. In the specific context of clinically diagnosing Post-Traumatic Stress Disorder (PTSD), subjective stressors, a high degree of similarity in presentation, the relatively easy access to information on how to fake PTSD, and financial incentives from a positive diagnosis in legal procedures or disability claims, is known to hinder accurate assessments. Here, we present an approach based on a generative model to flag malingering suspects rapidly, efficiently and accurately; essentially demonstrating a generative artificial intelligence (AI) framework that, paraphrasing Ekman⁴, “*can catch a liar*”.

Clinically, PTSD is an anxiety disorder that can develop after experiencing a traumatic event. In the United States, substantial disability compensation may be available for those with mental disorders, which while being an crucial resource for the truly afflicted, incentivizes malingering^{5,6}. Thus, faking PTSD to access medical treatment, commit insurance, personal injury and other frauds, or in an attempt to evade criminal liability and penalties^{7–10} is unfortunately not rare. While PTSD is a serious mental health condition associated with substance use disorder, mood disorder, anxiety disorder, personality disorder, increased morbidity, and possibly with increased mortality^{11,12}, false diagnoses can cause substantial financial drain^{13,14} to healthcare systems, divert crucial resources from where they are needed¹⁵, and interfere with study outcomes by introducing inaccuracies in clinical data¹⁶. Accurate disambiguation of true and fake PTSD is therefore of high importance, especially with sources suggesting that over 20% of personal injury cases, as well as 20% of the Veterans seeking combat compensation could be fabricating their condition^{5,17–19}. In addition to curbing financial fraud, ability to flag such incongruities can help address the myriad of complex factors that drive malingering behavior, including other mental health conditions, poverty, a lack of access to healthcare^{20,21} and a limited ability to seek help for mental health conditions.

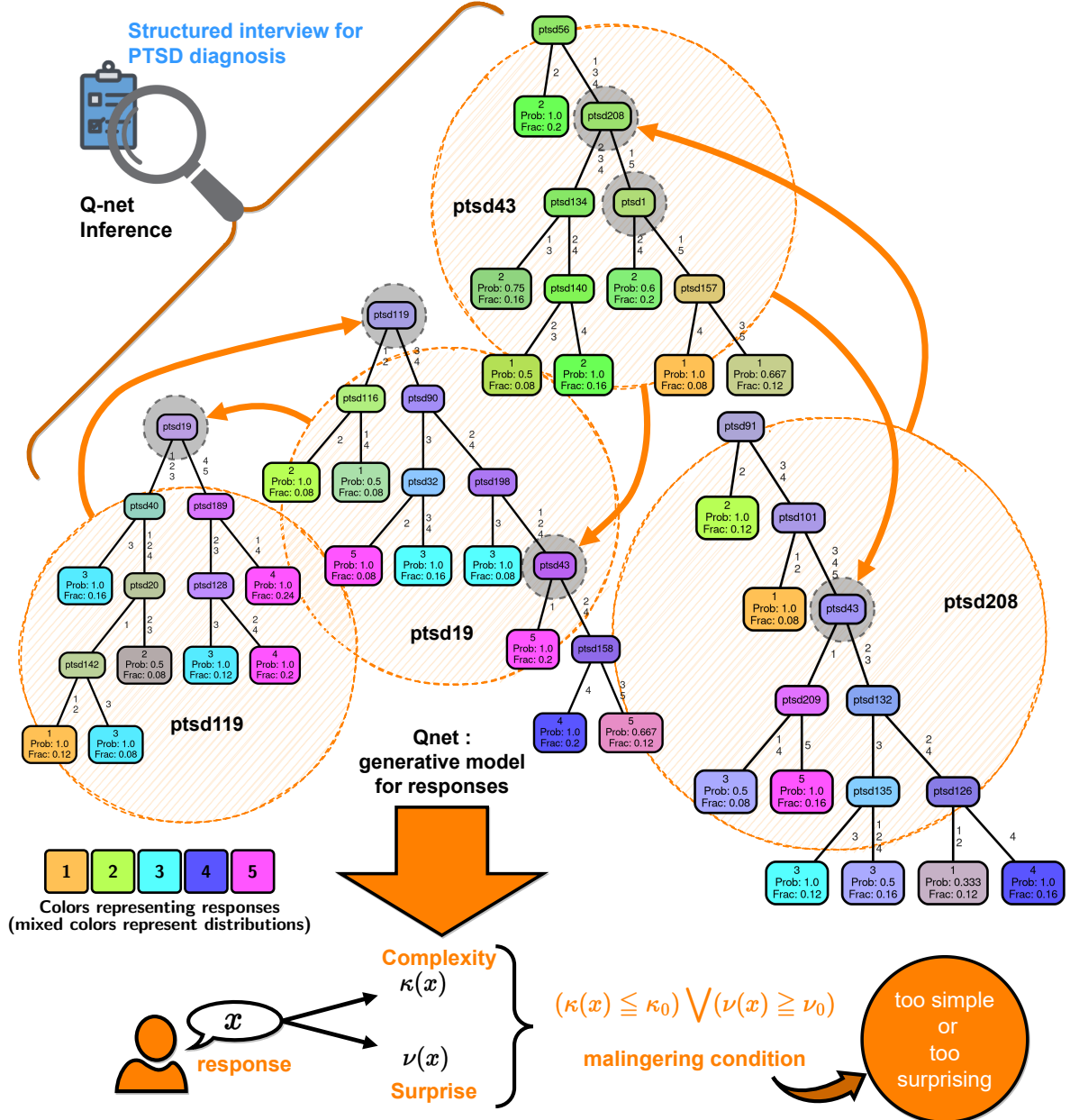


Fig. 1: Conceptual framework. Using a dataset of responses to a validated structured interview for PTSD diagnoses, along with physician-validated clinical diagnoses, we infer generative models for responses for PTSD patients. In our framework for detecting malingering, we flag responses as those which are highly “surprising” (defined as violating inferred cross-dependencies between individual response items) or are too simple (lacking complexity in the response patterns typical of non-malingering responses). The precise “malingering condition” shown above is validated from theoretical considerations as well as field experimental data.

Multiple standardized tests^{22,23} and validity assessment tools²⁴ have been explored for detecting malingering with limited success. These tools typically aim to incorporate patterns observed in diagnostic populations that might disambiguate faked symptoms from real ones, or ask similar or related questions multiple times to verify consistency. However, existing approaches do not target specific disorders, require expert interpretation, are almost always subjective, and by design are unlikely to be effective against a malingerer with psychiatric training (See Table 1). Other strategies that incorporate diverse physiological monitoring and linguistic analysis^{25–28} cannot be easily adopted in the present context of structured or semi-structured interviews.

In algorithm `VeRITAS` we leverage the fact that responses in a structured interview have emergent statistical dependencies arising from the nature of the questions themselves, and are modulated by the trait we are aiming to detect *e.g.*, PTSD pathology. These dependencies are too complex to easily mimic on-the-fly, even with knowledge about how such patients tend to respond *i.e.* with training in the mental health services. We find that `VeRITAS` offers a robust generalizable approach to verifying response validity, can be trained to target specific disorders, can be administered in less than 2 minutes, and require little or no subjective interpretation.

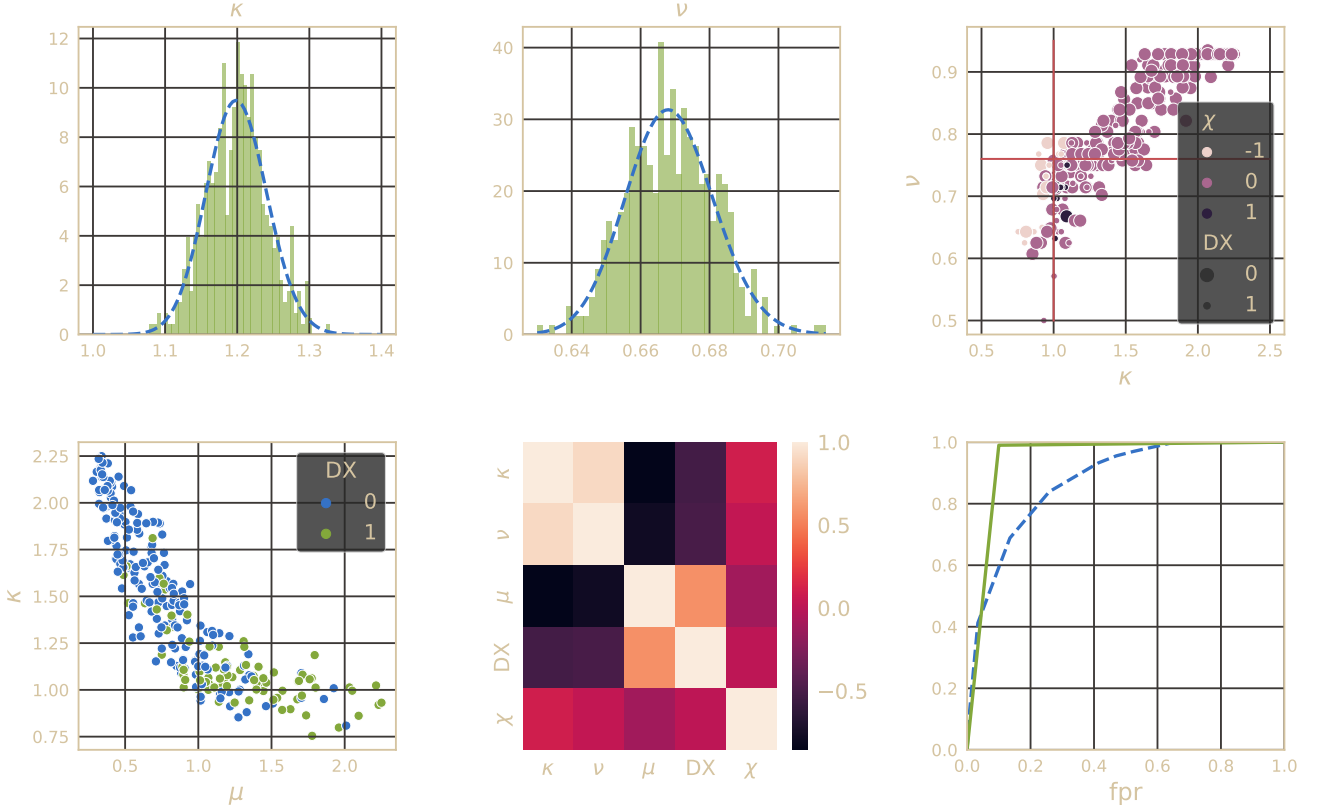


Fig. 2:

RESULTS

Participants and Data Sources

The interview data used in training in this study consists of $n = 304$ participants recruited at a Veterans Health Administration facility for an earlier study²⁹. Veterans were eligible if they were between the ages of 18 and 89 years and able to provide written informed consent. Once eligibility was determined by the study team, participants completed a PTSD-symptom questionnaire from the CAT-PTSD item bank^{29,30} (comprising 211 items) and provided responses on the PTSD Checklist for Diagnostic and Statistical Manual of Mental Disorders, Fifth Edition (DSM-5) (PCL-5). Participants in the training dataset were interviewed using the Clinician-Administered Scale for PTSD for DSM-5 (86 were diagnosed, 218 were deemed PTSD-negative).

Key Insight

Our key insight is that malingered responses have low likelihood of being generated while respecting the emergent cross-dependencies. The proposed framework enables verifying this condition efficiently without strong a priori assumptions of the nature and structure of these cross-dependencies. This condition can also be thought of as stipulating that malingered responses have high average “surprise”, *i.e.* deviate more on average from the model predictions of item responses. We learn a generative model for responses to structured diagnostic interviews that can discover n -way dependencies between the items. Finally, we aim to robustly flag malingering by checking if two specific response attributes κ, ν related to likelihood and surprise violate preset thresholds.

Q-net Inference and Cross-talk Modeling

Training and Validation

Performance

Comparison Against State of Art

DISCUSSION

Malingering remains a persistent problem, especially in psychiatric and criminal justice settings, with prevalence estimates ranging from 8 – 64%^{31–33}. Here we develop a novel system to address this gap for the need for further development of reliable, rigorous, and principled methods for detecting deception remains highly significant³⁴.

In this study, we introduce a novel algorithmic technique for identifying deception in structured interviews. This approach first infers a detailed model of the response patterns of previous respondents, then uses a measure of consistency known as dissonance to identify if future respondents' answers are indicative of deception.

Our model-based, data-driven approach is notably distinct from traditional methods of malingering detection, which are typically constructed for specific contexts³⁵. In the context of psychiatric diagnosis detection of faked symptoms is often based on standardized tests; some commonly used ones are summarized in Table 1. Tests such as the Structures Interview of Reported Symptoms (SIRS)²², Structured Inventory of Malingered Symptomatology (SIMS)²³, and the Minnesota Multiphasic Personality Inventory-2 and its associated validity scales²⁴ draw extensively on research of both legitimate psychiatric patients and malingerers, and employ a combination of detection strategies focused on symptoms that are either unlikely to be presented by genuine patients, or that tend to be amplified in malingerers³⁶. For detection of cognitive malingering, physiological, arousal-based approaches such as the Control Question Technique³⁷, and cognitive load-inducing approaches such as requiring subjects to perform a concurrent task (*e.g.* maintaining eye contact) and the Autobiographical Implicit Association Test (aIAT) have been suggested.

Accuracy rates for these research-driven state of the art detection techniques are typically in the range of 85-95% depending on the specific problem context³⁸⁻⁴⁰ (similarly high sensitivities have also been reported²³). However, among the limitations to these methods is the requirement of substantial domain research and/or subject matter expertise to develop. Thus while these methods perform well in the specific context for which they were developed, it may be challenging and/or expensive to broadly extend to detection in other contexts. In the case of cognitive malingering methods, vulnerabilities to clinician/administrator bias, high rates of false positives and false negatives, and vulnerability to coaching have been noted as possible limitations.

In contrast with these methods, our results demonstrate that it is possible to achieve similar (or perhaps improved) performance in detecting malingering by identifying structural differences in an individual's responses compared to a set of baseline responses. The basic insight underpinning our results is that consistent, genuine responses to interview questions are typically not independent. In general, there are dependence patterns between interview questions, which are often nontrivial and difficult or impossible to identify *a priori*. Given the entirety of responses an individual has provided in an ongoing interview, the response probabilities for a consistent response to the next question are governed by these dependence patterns. As these probabilistic constraints exist between all questions in the interview, a complex interaction network is induced over the high-dimensional response space. To infer the structure of this network, we use a novel architecture - the Q-net, a nonparametric generative model of the response network which reflects the n -way inter-question relations found amongst previous interview respondents.

From the model's approximation of the probabilistic dependencies present in the response network, we derive the dissonance as an intuitive, individually computable measure of response consistency. For all data sets in this study, mean dissonance values for random malingering and expert malingering are significantly higher than for actual responses, empirically demonstrating its utility for identifying such responses. As should be expected, we also found that the expert malingering responses were on average less dissonant than random malingering, showing that domain knowledge does allow individuals to respond in a more convincing manner. However, their failure to adequately account for detailed inter-question dependencies tends to make their responses significantly different from actual responses in a measurable, quantitative sense.

The classification results obtained further demonstrate the ability of this framework to identify subtle deviations in the response patterns of actual and malingering respondents. Due to the lower dissonance generated by expert malingerers, we found it more difficult to identify such responses; however, in practice it appears unlikely that even the most knowledgeable domain experts would have the requisite quantitative knowledge of response distributions of the target phenotype to be able to respond in the manner of our idealized expert malingerers. Thus in real-world usage where true malingerers are likely a mixture of these two idealizations, we would expect performance to be bounded above and below by the empirical performance obtained for the random and expert scenarios.

In general, data-driven methods for detection of malingering are not as ubiquitous as in other domains; however, the increasing availability of modern sensory data collection has allowed for implementation of machine learning algorithms and computer vision techniques^{41,42} for detecting deception based on patterns in language, facial expressions, body movements, and eye movements^{43,44}. These methods have been shown to outperform traditional lie detection methods in some cases, obtaining accuracy rates of over 90% for lie detection from patterns in facial expressions and body movements^{4,43} and 60-80% using analysis of gaze direction and pupil dilation⁴⁴. However, in many instances it is not feasible to collect detailed sensory data, possibly limiting the general applicability of these kinds of methods. By exclusively utilizing response data, our dissonance-based approach is not subject to such limitations.

In contexts such as these, where malingerers seek to obtain a target diagnosis, training data from past respondents often includes a clinical phenotype for each respondent. However, this is not always available (as was observed for the CCHHS data set). Nonetheless, in the absence of the underlying diagnostic phenotypes,

our approach still appears to perform well under a scenario in which expert malingerers are also unsure of ground-truth diagnoses.

In other types of structured interviews, such as internet surveys and opinion polls, built-in data validity checks are frequently used to attempt to proactively identify disingenuous respondents. Among such checks, unusually fast response times have been observed to be indicative of poor data quality^{45,46}, demonstrating a lack of attention by interview respondents⁴⁷. Our finding of a negative correlation between dissonance and response time adds some support to this notion. However, it has been noted that duration-based validity checks may be vulnerable to a high false negative rate⁴⁸, suggesting a need for most robust measures of assessing validity.

METHODS

1. DEFINITIONS & NOTATION

Definition 1 (Survey). *A survey for the purpose of this work is a structured interview, consisting of a finite number of questions (items) posed to a set of participants, with these items drawn from a finite item bank, and whose responses must be one from a pre-specified set of choices, e.g., the Likert scale, with missing values for the responses allowed.*

Definition 2 (Response vector). *A response vector is the set of responses to a survey from a single participant, typically assuming that not all items are posed, and allows for the possibility that some responses are missing.*

A Q-net, as described here, is a model of the response dependency structure for questions (items) posed to participants in a survey. The Q-net explicitly estimates individual conditional distributions of each item response, which collectively serve as a model of the full joint distribution of the responses.

Definition 3 (Q-net). *Let $X \sim P$ be an n -dimensional discrete random vector supported on a finite set Σ and following distribution P , i.e.*

$$X = (X_1, \dots, X_n) \sim P, \quad \text{supp}(X) = \Sigma = \prod_{i=1}^n \Sigma_i \quad \text{with } |\Sigma| < \infty.$$

For $i = 1, \dots, n$, let $P_i := P(X_i | X_j = x_j \text{ for } j \neq i)$ denote the conditional distribution of X_i given the values of the other components of X . Finally, for each $i = 1, \dots, n$, let Φ_i^P denote an estimate of the distribution P_i . Then the set $\Phi^P := \{\Phi_i^P\}_{i=1}^n$ is called a Quasinet (Q-net). Identifying the true distribution P as the one describing the joint statistics of the responses from a survey with n items, we also refer to Φ^P as the Q-net for the survey P .

When P is clear from context, we may omit the superscript and simply write $\Phi = \{\Phi_i\}$ to denote the Q-net. The motivation for Definition 3 is that the collection of all estimators $\Phi = \{\Phi_i\}$ contained in a Q-net represents the set of all inferred dependencies from the observed ecosystem. While the definition allows for arbitrary method of algorithm to construct the estimators Φ_i , the utility of a Q-net clearly depends primarily on the properties of the Φ_i . In this study, we aim to minimize the set of a priori assumptions on the overall model structure to allow the complex dependencies present in P to emerge. To that end, throughout this work all Q-nets are computed using conditional inference trees⁴⁹ (a variant of classification and regression trees) to compute each Φ_i . In general each Q-net component Φ_i is computed independently from the other Φ_j , which allows a network structure to emerge amongst these estimators.

An important quantity for an inferred Q-net is the persistence function ω_x .

Definition 4 (Persistence Function). *Given a survey P inducing the Q-net Φ^P and a response vector $x = (x_1, \dots, x_n)$, the persistence ω_x of x in the population modeled by the Q-net:*

$$\omega_x^P := \Pr(x \in P) = \prod_{i=1}^n \Phi_i^P(X_i = x_i | X_j = x_j, j \neq i) \quad (1)$$

The persistence function ω_x^P , as the name suggests, is the probability that x persists, i.e., $\Pr(x \rightarrow x)$ for the population modeled by the Q-net P , with $1 - \omega_x^P$ being the probability that x is altered by a random perturbation.

We will show that if for two inferred Q-net models P, Q , we have $\omega_x^P \geq \omega_x^Q$, then it is more likely that model P generated x . This is an important result that justifies the definition of the score parameter in Defn. 6.

The Q-net allows us to rigorously compute bounds on the probability of a spontaneous change from one response vector to another, induced by spontaneous chance variations (Fig ??). Not all perturbations in a vector are either likely or contextually meaningful. With an exponentially exploding number of possibilities in which a vector over a large set of items can vary, it is computationally intractable to directly model all possible dependencies; nevertheless, we can constrain the possibilities using the patterns we uncover via the Q-net construction. A key

piece of this approach is to design an intrinsic distance between response vectors, which is reflective of this underlying dependency structure.

Definition 5 (q-distance). Let $\Phi^P = \{\Phi_i^P\}_{i=1}^n$ and $\Phi^Q = \{\Phi_i^Q\}_{i=1}^n$ denote Q-nets on populations P and Q , and suppose $\mathbf{x} = (x_1, \dots, x_n)$ and $\mathbf{y} = (y_1, \dots, y_n)$ are samples of $X \sim P$ and $Y \sim Q$ respectively. Then the q-distance $\theta_{P,Q}(\mathbf{x}, \mathbf{y})$ between \mathbf{x} and \mathbf{y} is

$$\theta_{P,Q}(\mathbf{x}, \mathbf{y}) := \frac{1}{n} \sum_{i=1}^n \left[\mathbb{J}^{\frac{1}{2}} \left(\Phi_i^P(X_i | X_j = x_j, j \neq i) \parallel \Phi_i^Q(Y_i | Y_j = y_j, j \neq i) \right) \right]$$

where \mathbb{J} denotes the Jensen-Shannon divergence⁵⁰.

For brevity, we may write simply θ (dropping the suffixes) if the populations are clear from context. Since the Jensen-Shannon distance \mathbb{J} is a legitimate metric⁵¹ on the set of probability distributions (unlike KL-divergence), θ inherits nonnegativity, symmetry, and respects the triangle inequality; it follows that q-distance is a (pseudo)-metric on Σ . Note that, being a pseudo-metric implies that we may have $\theta(\mathbf{x}, \mathbf{y}) = 0$ for $\mathbf{x} \neq \mathbf{y}$, i.e. distinct vectors can induce the same distributions over each index, and thus have zero distance. This is in fact desirable, since we do not want our distance to be sensitive to changes that are not meaningful. The intuition is that not all variations are equally important or likely. Moreover, we show in Theorem 1 that the log-likelihood of a vector \mathbf{x} transitioning to \mathbf{y} scales with $\theta(\mathbf{x}, \mathbf{y})$, allowing us to directly estimate the probability of spontaneous (or sequential) jumps between abundance profiles.

Theorem 1 (Probability Bound). Given a vector \mathbf{x} of length n from P that transitions to \mathbf{y} from Q , we have the following bounds at significance level α .

$$\omega_y e^{\frac{\sqrt{8}N^2}{1-\alpha}\theta(\mathbf{x}, \mathbf{y})} \geq \Pr(\mathbf{x} \rightarrow \mathbf{y}) \geq \omega_y e^{-\frac{\sqrt{8}N^2}{1-\alpha}\theta(\mathbf{x}, \mathbf{y})} \quad (2)$$

where ω_y is the persistence of \mathbf{y} (Def. 4), and $\theta(\mathbf{x}, \mathbf{y})$ is the q-distance between \mathbf{x}, \mathbf{y} (Def. 5).

Proof. See later in Section 3. □

Theorem 1 gives theoretical backing to the claim that samples generated by the Q-net indeed reflect likely perturbation possibilities from the current state. Thus we can use the Q-net to draw contextually realistic samples that respect the cross dependencies and reduce surprise (that is, the Q-net-inferred conditional distributions can be used to generate approximate samples from the population P). This has several implications, such as the ability to easily handle missing/incomplete data.

Remark 1 (Neighborhood Structure). It follows from Th. 1 that we have for some constant C ,

$$\ln \left| \frac{\Pr(\mathbf{x} \rightarrow \mathbf{y})}{\Pr(\mathbf{y} \rightarrow \mathbf{x})} \right| \leq C\theta(\mathbf{x}, \mathbf{y}) \quad (3)$$

implying for all response vectors \mathbf{y} within a small neighborhood of \mathbf{x} (small in metric θ), we have:

$$\Pr(\mathbf{y} \rightarrow \mathbf{x}) \approx \Pr(\mathbf{x} \rightarrow \mathbf{x}) \quad (4)$$

which reveals an important special structure on local neighborhoods.

2. VERITAS ANALYSIS

Definition 6 (Algorithm VERITAS Parameters). We introduce three parameters referred to as the complexity, surprise and score parameters (κ, ν, μ respectively) for a given response vector \mathbf{x} :

$$\text{complexity: } \kappa \triangleq -\frac{1}{|\mathbf{x}|} \ln \Pr(\mathbf{x} \rightarrow \mathbf{x} | M^+) = -\frac{\ln \omega_x^{M^+}}{|\mathbf{x}|} \quad (5)$$

$$\text{surprise: } \nu \triangleq \mathbf{E}_i \left(1 - \Phi_i^{M^+}(\mathbf{x}_{-i}) | \mathbf{x}_i \right) \quad (6)$$

$$\text{score: } \mu \triangleq \frac{\ln \Pr(\mathbf{x} \rightarrow \mathbf{x} | M^+)}{\ln \Pr(\mathbf{x} \rightarrow \mathbf{x} | M^0)} = \frac{\ln \omega_x^{M^+}}{\ln \omega_x^{M^0}} \quad (7)$$

where M^+ indicates the sub-population exhibiting a particular trait of interest e.g. a mental health disorder such as PTSD, and M^0 is the control sub-population where this trait is absent.

Interpretation of the Parameters

We

Definition 7 (Malingering property). A response vector \mathbf{x} is defined to have the malingering property if:

$$(\mu(\mathbf{x}) \geq \mu_0) \bigwedge \left((\kappa(\mathbf{x}) \leq \kappa_0) \vee (\nu(\mathbf{x}) \geq \nu_0) \right) \quad (8)$$

Set of malingering responses is denoted as \mathcal{M} .

The decision thresholds κ_0, ν_0, μ_0 are inferred from survey data.

Lemma 1 (Complexity). *For a survey with n items, and assuming L to be the number of possible responses to each item, the unconditional probability of a response vector x occurring among all feasible responses is bounded above by $(e^\kappa/L)^n$, where $\kappa(x)$ is the complexity parameter for response x .*

Proof. Let $\kappa(x) \leq \kappa'$. From Def. 6, we have for a response vector x ,

$$-\frac{1}{n} \ln \omega_x \leq \kappa' \Rightarrow \omega_x \geq e^{-n\kappa'} \quad (9)$$

Summing on both sides over all responses x with $\kappa(x) \leq \kappa'$ (assume there are N_x such sequences), we have:

$$1 \geq \sum_x \omega_x \geq \sum_x e^{-n\kappa'} \quad (10)$$

where the first inequality follows from observing that responses very close to x in the q-distance metric have a specific structure, namely $\omega_x \approx Pr(y \rightarrow x)$ (See Remark 1) and responses further away have smaller jump probabilities, which then implies:

$$N_x \sum_x e^{-n\kappa'} \leq 1 \Rightarrow N_x \leq e^{n\kappa'} \quad (11)$$

The result then follows from noting that the complete set of possible responses has the size L^n . \square

Lemma 1 justifies why a low value of κ implies the possibility of an un-natural response, because the odds of generating such a response is remarkably small.

Lemma 2 (Surprise). *For any response vector x , we have:*

$$\nu(x) \leq 1 - e^{-\kappa(x)} \quad (12)$$

Proof. Denoting $\Phi_i(x_{-i})|_{x_i}$ as a_i , we note that $\omega_x^{1/n}$ is the geometric mean of the vector of a_i s, while $\mathbf{E}_i(\Phi_i(x_{-i})|_{x_i})$ is the arithmetic mean of the same vector, which then completes the proof by noting:

$$-\mathbf{E}_i(\Phi_i(x_{-i})|_{x_i}) \leq -\omega_x^{1/n} \Rightarrow \nu(x) \leq 1 - \omega_x^{1/n} \quad (13)$$

\square

Why the Defined Property Identifies Malingering

Lemma 2 indicates that the requirement of an upper bound on the surprise and a lower bound on the complexity are both aiming to flag responses which are unlikely to appear when the data (responses) are being generated by the underlying process corresponding to the phenotype of interest (PTSD). When such unlikely responses do appear nevertheless, it is likely that they are not being generated by the correct underlying process. Can it correspond to another “natural” process, *e.g.* a different mental disorder that is not intentional fabrication? If the score is higher than

Note that the remaining condition $\mu(x) \geq \mu_0$ is a diagnosis criterion for the trait of interest (M^+), and may be replaced with a different condition if available for identifying participants with the M^+ trait. This particular form follows from a straightforward Bayesian argument on estimating the posterior.

3. PROOF OF THEOREM 1

Theorem 2 (Probability bound). *Given a sequence x of length N that transitions to a strain $y \in Q$, we have the following bounds at significance level α .*

$$\omega_y^Q e^{\frac{\sqrt{8}N^2}{1-\alpha}\theta(x,y)} \geq Pr(x \rightarrow y) \geq \omega_y^Q e^{-\frac{\sqrt{8}N^2}{1-\alpha}\theta(x,y)} \quad (14)$$

where ω_y^Q is the membership probability of strain y in the target population Q (See Def. ??), and $\theta(x, y)$ is the q-distance between x, y (See Def. 5).

Proof. Using Sanov’s theorem⁵⁰ on large deviations, we conclude that the probability of spontaneous jump from strain $x \in P$ to strain $y \in Q$, with the possibility $P \neq Q$, is given by:

$$Pr(x \rightarrow y) = \prod_{i=1}^N (\Phi_i^P(x_{-i})|_{y_i}) \quad (15)$$

Writing the factors on the right hand side as:

$$\Phi_i^P(x_{-i})|_{y_i} = \Phi_i^Q(y_{-i})|_{y_i} \left(\frac{\Phi_i^P(x_{-i})|_{y_i}}{\Phi_i^Q(y_{-i})|_{y_i}} \right) \quad (16)$$

we note that $\Phi_i^P(x_{-i})$, $\Phi_i^Q(y_{-i})$ are distributions on the same index i , and hence:

$$|\Phi_i^P(x_{-i})_{y_i} - \Phi_i^Q(y_{-i})_{y_i}| \leq \sum_{y_i \in \Sigma_i} |\Phi_i^P(x_{-i})_{y_i} - \Phi_i^Q(y_{-i})_{y_i}| \quad (17)$$

Using a standard refinement of Pinsker's inequality⁵², and the relationship of Jensen-Shannon divergence with total variation, we get:

$$\theta_i \geq \frac{1}{8} |\Phi_i^P(x_{-i})_{y_i} - \Phi_i^Q(y_{-i})_{y_i}|^2 \Rightarrow \left| 1 - \frac{\Phi_i^Q(y_{-i})_{y_i}}{\Phi_i^P(x_{-i})_{y_i}} \right| \leq \frac{1}{a_0} \sqrt{8\theta_i} \quad (18)$$

where a_0 is the smallest non-zero probability value of generating the entry at any index. We will see that this parameter is related to statistical significance of our bounds. First, we can formulate a lower bound as follows:

$$\log \left(\prod_{i=1}^N \frac{\Phi_i^P(x_{-i})_{y_i}}{\Phi_i^Q(y_{-i})_{y_i}} \right) = \sum_i \log \left(\frac{\Phi_i^P(x_{-i})_{y_i}}{\Phi_i^Q(y_{-i})_{y_i}} \right) \geq \sum_i \left(1 - \frac{\Phi_i^Q(y_{-i})_{y_i}}{\Phi_i^P(x_{-i})_{y_i}} \right) \geq \frac{\sqrt{8}}{a_0} \sum_i \theta_i^{1/2} = -\frac{\sqrt{8}N}{a_0} \theta \quad (19)$$

Similarly, the upper bound may be derived as:

$$\log \left(\prod_{i=1}^N \frac{\Phi_i^P(x_{-i})_{y_i}}{\Phi_i^Q(y_{-i})_{y_i}} \right) = \sum_i \log \left(\frac{\Phi_i^P(x_{-i})_{y_i}}{\Phi_i^Q(y_{-i})_{y_i}} \right) \leq \sum_i \left(\frac{\Phi_i^Q(y_{-i})_{y_i}}{\Phi_i^P(x_{-i})_{y_i}} - 1 \right) \leq \frac{\sqrt{8}N}{a_0} \theta \quad (20)$$

Combining Eqs. 19 and 20, we conclude:

$$\omega_y^Q e^{\frac{\sqrt{8}N}{a_0} \theta} \geq Pr(x \rightarrow y) \geq \omega_y^Q e^{-\frac{\sqrt{8}N}{a_0} \theta} \quad (21)$$

Now, interpreting a_0 as the probability of generating an unlikely event below our desired threshold (*i.e.* a "failure"), we note that the probability of generating at least one such event is given by $1 - (1 - a_0)^N$. Hence if α is the pre-specified significance level, we have for $N \gg 1$:

$$a_0 \approx (1 - \alpha)/N \quad (22)$$

Hence, we conclude, that at significance level $\geq \alpha$, we have the bounds:

$$\omega_y^Q e^{\frac{\sqrt{8}N^2}{1-\alpha} \theta} \geq Pr(x \rightarrow y) \geq \omega_y^Q e^{-\frac{\sqrt{8}N^2}{1-\alpha} \theta} \quad (23)$$

□

TABLE 1: Summary of malingering/deception detection methods.

Approach/Tool	Description	Noted Accuracy/Performance
Structures Interview of Reported Symptoms (SIRS) ²²	An interview-based measure with multiple detection strategies.	No specific accuracy rate mentioned, but noted as a robust instrument.
Structured Inventory of Malingered Symptomatology (SIMS) ²³	A paper-and-pencil screening device for detecting malingering with a sensitivity for detecting malingering of 95.6% in a study with college students.	95.6% sensitivity in a specific study context.
Minnesota Multiphasic Personality Inventory-2 (MMPI-2) ²⁴	A self-report measure assessing personality and psychopathology. Certain validity scales were developed to uncover malingering.	Not specified, but noted flaws and potential for false positives.
Millon Clinical Multiaxial Inventory MCMI-III	A self-report scale focusing on personality disorders.	No specific accuracy rate mentioned.
Miller Forensic Assessment of Symptoms (M-FAST)	A brief screening measure for malingered mental illness in forensic settings.	Noted issues with low internal consistency on some scales.
Human Lie Detectors	Not a specific tool, but the general method of humans attempting to discern lies from truth.	People are generally poor lie detectors.
Arousal-Based Approaches (like the Polygraph)	Techniques that rely on physiological responses.	Criticized for poor validity and high rate of false positives.

Cognitive Load-Inducing Approaches	Techniques that view deception as a cognitive act that generally imposes greater cognitive load on respondents than honesty does.	No specific accuracy rate mentioned.
Autobiographical Implicit Association Test (aIAT) ³⁸	Designed to determine whether respondents possess actual autobiographical memories.	91% accuracy rate in identifying genuine autobiographical memories.
Timed Antagonistic Response Alethiometer (TARA) ³⁹	A computer-administered, response time-based method of lie detection.	85% accuracy rate.
Detecting Faked Identities With Unexpected Questions and Mouse Movements	Technique using computer mouse movements in conjunction with unexpected questions to uncover faked identities ⁴⁰	95% accuracy rate.
Time-Restricted Integrity Confirmation (TRI-Con)	A cognitive load-inducing technique with potential to uncover different kinds of deception including malingering.	Up to 89% accuracy rate.
Activation-Decision-Construction-Action Theory (ADCAT)	A theory of high-stakes deception.	No specific accuracy rate mentioned, but this is more of a theoretical foundation rather than a specific tool or method.

REFERENCES

- [1] Ali, S., Jabeen, S. & Alam, F. Multimodal approach to identifying malingered posttraumatic stress disorder: A review. *Innovations in Clinical Neuroscience* **12**, 12 (2015).
- [2] Rogers, R. An overview of malingering and its assessment. *Psychiatric Clinics of North America* **20**, 15–27 (1997).
- [3] Rogers, R. E. *Clinical assessment of malingering and deception* (Guilford Press, 2008).
- [4] Ekman, P. & O'Sullivan, M. Who can catch a liar? *American Psychologist* **46**, 913–920 (1991).
- [5] Frueh, B., Grubaugh, A., Elhai, J. & Buckley, T. Us department of veterans affairs disability policies for posttraumatic stress disorder: administrative trends and implications for treatment, rehabilitation, and research. *Am J Public Health* **97**, 2143–5 (2007).
- [6] Taylor, S., Frueh, B. C. & Asmundson, G. J. G. Detection and management of malingering in people presenting for treatment of posttraumatic stress disorder: Methods, obstacles, and recommendations. *Journal of Anxiety Disorders* **21**, 22–41 (2007).
- [7] Guriel, J. & Fremouw, W. Assessing malingered posttraumatic stress disorder: A critical review. *Clinical Psychology Review* **23**, 881–904 (2003).
- [8] Salloway, S., Southwick, S. & Sadowsky, M. Opiate withdrawal presenting as posttraumatic stress disorder. *Hospital and Community Psychiatry* **41**, 666–667 (1990).
- [9] Resnick, P. J., West, S. & Payne, J. W. Malingering of posttraumatic disorders. In Rogers, R. (ed.) *Clinical assessment of malingering and deception*, 109–127 (Guilford Press, 2008), 3 edn.
- [10] Burkett, B. G. & Whitley, G. *Stolen valor: How the Vietnam generation was robbed of its heroes and history* (Verity Press, 1998).
- [11] Goldstein, R. B. *et al.* The epidemiology of dsm-5 posttraumatic stress disorder in the united states: results from the national epidemiologic survey on alcohol and related conditions-iii. *Social psychiatry and psychiatric epidemiology* **51**, 1137–1148 (2016).
- [12] Schnurr, P. P., Lunney, C. A., Bovin, M. J. & Marx, B. P. Posttraumatic stress disorder and quality of life: Extension of findings to veterans of the wars in iraq and afghanistan. *Clinical psychology review* **29**, 727–735 (2009).
- [13] LoPiccolo, C., Goodkin, K. & Baldewicz, T. Current issues in the diagnosis and management of malingering. *Ann Med* **31**, 166–174 (1999).
- [14] Oboler, S. Disability evaluations under the department of veterans affairs. In Rondinelli, R. & Katz, R. (eds.) *Impairment Rating and Disability Evaluation*, 187–217 (W. B. Saunders, Philadelphia, PA, 2000).
- [15] Taylor, S. *Clinician's Guide to Treating PTSD: A Cognitive-Behavioral Approach* (Guilford Press, New York, 2006).
- [16] Rosen, G. Dsm's cautionary guideline to rule out malingering can protect the ptsd data base. *J Anxiety Disorders* **20**, 530–535 (2006).
- [17] Marx, B. & Holowka, D. Ptsd disability assessment. *PTSD Res Q* **22**, 1–6 (2011).

- [18] Rogers, R., Sewell, K. & Goldstein, A. Explanatory models of malingering: a prototypical analysis. *Law & Hum Behav* **18**, 543–52 (1994).
- [19] Lees-Haley, P. Mmpi-2 base rates for 492 personal injury plaintiffs: implications and challenges for forensic assessment. *J Clin Psychol* **53**, 745–55 (1997).
- [20] Park, L., Costello, S., Li, J., Lee, R. & Jacobson, K. C. Race, health, and socioeconomic disparities associated with malingering in psychiatric patients at an urban emergency department. *General Hospital Psychiatry* **71**, 121–127 (2021).
- [21] Muntaner, C., Eaton, W. W., Miech, R. & O'campo, P. Socioeconomic position and major mental disorders. *Epidemiologic reviews* **26**, 53–62 (2004).
- [22] Wong, S. & O'Sullivan, M. The structured interview of reported symptoms (sirs): An overview. *Assessment* **12**, 289–307 (2005).
- [23] Smith, G. P. & Burger, G. K. Detection of malingering: validation of the structured inventory of malingered symptomatology (sims). *Journal of the American Academy of Psychiatry and the Law Online* **25**, 183–189 (1997).
- [24] Ben-Porath, Y. S. *Interpreting the mmpi-2-rf* (U of Minnesota Press, 2012).
- [25] Ekman, P. & O'Sullivan, M. Who can catch a liar? *American Psychologist* **46**, 913 (1991).
- [26] Mihalcea, R. & Strapparava, C. The lie detector: Explorations in the automatic recognition of deceptive language. In *ACL-IJCNLP 2009* (2009).
- [27] Burgoon, J. K., Blair, J. P. & Strom, R. E. Cognitive biases and nonverbal cue availability in detecting deception. *Human Communication Research* **34**, 572–599 (2008).
- [28] Zhou, L., Burgoon, J. K., Nunamaker Jr, J. F. & Twitchell, D. Automating linguistics-based cues for detecting deception in text-based asynchronous computer-mediated communications. *Group Decision and Negotiation* **13**, 81–106 (2004).
- [29] Brenner, L. A. *et al.* Development and validation of computerized adaptive assessment tools for the measurement of posttraumatic stress disorder among us military veterans. *JAMA Network Open* **4**, e2115707–e2115707 (2021).
- [30] Gibbons, R. D. *et al.* Development of a computerized adaptive test for depression. *Archives of general psychiatry* **69**, 1104–1112 (2012).
- [31] McDermott, B., Dualan, I. & Scott, C. Malingering in the correctional system: does incentive affect prevalence? *Int'l J L & Psychiatry* **36**, 287–92 (2013).
- [32] Schmidt, T., Krüger, M. & Ullmann, U. Base rate of probable malingering and its indicators in the assessment of mental disorders-retrospective analysis of a sample of forensic psychological evaluations. *Die Rehabilitation* **59**, 231–236 (2020).
- [33] Matto, M., McNiel, D. E. & Binder, R. L. A systematic approach to the detection of false ptsd. *The journal of the American Academy of Psychiatry and the Law* **47**, 325–334 (2019).
- [34] DePaulo, B. M., Lindsay, J. J., Malone, B. E., Muhlenbruck, L. & Charlton, K. Cues to deception. *Psychological Bulletin* **129**, 74–118 (2003).
- [35] Walczyk, J. J., Sewell, N. & DiBenedetto, M. B. A review of approaches to detecting malingering in forensic contexts and promising cognitive load-inducing lie detection techniques. *Frontiers in psychiatry* **9**, 700 (2018).
- [36] Rogers, R. & Correa, A. A. Determinations of malingering: Evolution from case-based methods to detection strategies. *Psychiatry, Psychology and Law* **15**, 213–223 (2008).
- [37] Vrij, A. & Mann, S. Who killed my relative? police officers' ability to detect real-life high-stake lies. *Psychology, Crime & Law* **7**, 119–132 (2001).
- [38] Sartori, G., Agosta, S., Zogmaister, C., Ferrara, S. D. & Castiello, U. How to accurately detect autobiographical events. *Psychological science* **19**, 772–780 (2008).
- [39] Gregg, A. P. When vying reveals lying: The timed antagonistic response alethiometer. *Applied Cognitive Psychology: The Official Journal of the Society for Applied Research in Memory and Cognition* **21**, 621–647 (2007).
- [40] Monaro, M., Gamberini, L. & Sartori, G. The detection of faked identity using unexpected questions and mouse dynamics. *PloS one* **12**, e0177851 (2017).
- [41] Chan, E. Y., Danker-Hopfe, H. & Riemann, D. A machine learning approach to polysomnographic sleep study for lie detection. *IEEE Journal of Biomedical and Health Informatics* **22**, 716–722 (2018).
- [42] Wang, X., Yang, Y., Zhang, J., Wang, L. & Liu, S. A machine learning approach to lie detection based on eye movement and gesture analysis. *IEEE Transactions on Human-Machine Systems* **50**, 266–273 (2020).
- [43] Yu, H., Chen, Y. & Fu, X. Deep learning for deception detection: a comprehensive review. *IEEE Transactions on Information Forensics and Security* **13**, 2615–2629 (2018).
- [44] Ducharme, J. J., Langleben, D. D. & Verma, R. Eye tracking in the detection of deception: a review of the literature. *Psychology, Crime & Law* **23**, 875–895 (2017).
- [45] Malhotra, N. Completion time and response order effects in web surveys. *Public opinion quarterly* **72**, 914–934 (2008).
- [46] Callegaro, M., Yang, Y., Bhola, D. S., Dillman, D. A. & Chin, T.-Y. Response latency as an indicator of optimizing in online questionnaires. *Bulletin of Sociological Methodology/Bulletin de Méthodologie*

Sociologique **103**, 5–25 (2009).

- [47] Greszki, R., Meyer, M. & Schoen, H. Exploring the effects of removing “too fast” responses and respondents from web surveys. *Public Opinion Quarterly* **79**, 471–503 (2015).
- [48] Kennedy, C. *et al.* Assessing the risks to online polls from bogus respondents. *Pew Research Center* **18** (2020).
- [49] Sarda-Espinosa, A., Subbiah, S. & Bartz-Beielstein, T. Conditional inference trees for knowledge extraction from motor health condition data. *Engineering Applications of Artificial Intelligence* **62**, 26–37 (2017).
- [50] Cover, T. M. & Thomas, J. A. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)* (Wiley-Interscience, New York, NY, USA, 2006).
- [51] Fuglede, B. & Topsøe, F. Jensen-shannon divergence and hilbert space embedding. In *International Symposium on Information Theory, 2004. ISIT 2004. Proceedings.*, 31 (IEEE, 2004).
- [52] Fedotov, A. A., Harremoës, P. & Topsøe, F. Refinements of pinsker’s inequality. *IEEE Transactions on Information Theory* **49**, 1491–1498 (2003).